# INTRODUCTION TO INTEL® VTUNE™ PROFILER & INTEL® ADVISOR

Matthew Cordery

Application Engineer

matthew.cordery@intel.com

# Notices and Disclaimers

**DISTRIBUTION STATEMENT: None Required**

# Get the tools

- Advisor and VTune are now part of the Intel® oneAPI Base Toolkit

  - Download entire toolkit

  - Download just VTune and Advisor (customizable installation)

  - https://software.intel.com/content/www/us/en/develop/tools/oneapi/base-toolkit/download.html

  - Available for Windows, Linux, MacOS (view Linux results)

# Agenda

- Advisor – Intel's vectorization and optimization tool

  - CPU Optimization

    - Roofline

  - GPU Offloading

    - Offload Advisor

    - Roofline

- VTune – Intel's performance metric investigation tool

  - CPU capabilities

  - GPU metrics

# Tuning at Multiple Hardware Levels

Exploiting all features of modern processors requires good use of the available resources

- Core
  - Vectorization is critical with 512bit FMA vector units (32 DP ops/cycle)
  - Cache use needed to feed vector units
- Socket
  - Using all cores in a processor requires parallelization (MPI*, OMP*,CUDA*,OPENCL*,SYCL*,DPC++ ... )
  - Using coherent, shared socket caches
- Node
  - Minimize remote memory access (control memory affinity)
  - Minimize resource sharing (tune local memory access, disk IO and network traffic)

# ADVISOR: NBODY DEMONSTRATION

The naïve code that could

# N-body code

- Dr. Fabio Baruffa (original): https://github.com/fbaru-dev/nbody-demo

- Paulius Velesko (includes gpu): https://github.com/pvelesko/nbody-demo.git
  - Basically, the code in this demo

# Nbody gravity simulation

Consider a distribution of $n$ point masses located at $r_i$ with masses $m_i$ and velocities and accelerations $v_i$ and $a_i$, respectively

We want to calculate the position of the particles after a certain time interval using Newton's law of gravity.

```cpp
struct Particle
{
 public:
   Particle() { init();}
   void init()
   {
     pos[0] = 0.; pos[1] = 0.; pos[2] = 0.;
     vel[0] = 0.; vel[1] = 0.; vel[2] = 0.;
     acc[0] = 0.; acc[1] = 0.; acc[2] = 0.;
     mass   = 0.;
   }
   real_type pos[3];
   real_type vel[3];
   real_type acc[3];
   real_type mass;
};
```

```cpp
for (i = 0; i < n; i++){          // update acceleration
   for (j = 0; j < n; j++){
      real_type distance, dx, dy, dz;
      real_type distanceSqr = 0.0;
      real_type distanceInv = 0.0;

     dx = particles[j].pos[0] - particles[i].pos[0];
     …

     distanceSqr = dx*dx + dy*dy + dz*dz + softeningSquared;
    distanceInv = 1.0 / sqrt(distanceSqr);

     particles[i].acc[0] += dx * G * particles[j].mass *
                    distanceInv * distanceInv * distanceInv;
     particles[i].acc[1] += …
     particles[i].acc[2] += …
```
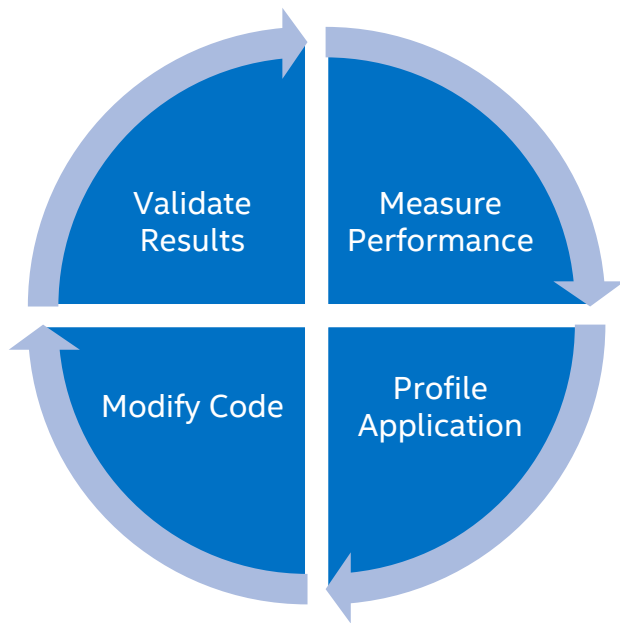
# The Basic Tuning Cycle



Validate Results → Measure Performance → Profile Application → Modify Code

Infinite cycle only broken by external constraints (time, papers, releases ... )

Procedures for measuring performance and validating results are critical

**Automation** and **environment** control are key for **consistency**

Where do I start?

# Version Optimizations

- Ver0
  - Initial implementation

- Ver1
  - Vectorized with compiler flags (march/mtune)

- Ver2
  - Use only floats

- Ver3/4
  - AoS -> SoA + SIMD Reduce

- Ver 7
  - OpenMP with data alignment

# INTEL® ADVISOR

Vectorization and Static Analysis

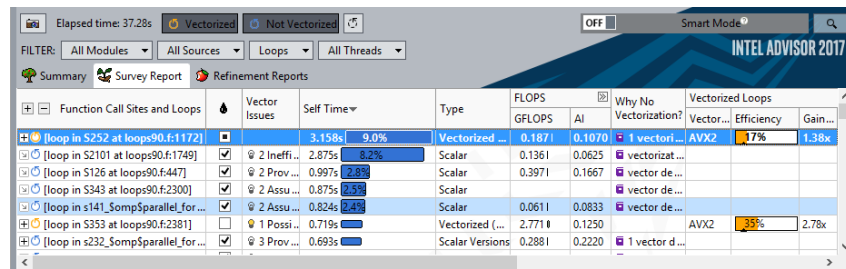https://www.alcf.anl.gov/user-guides/advixe-cl-xc40

# Intel® Advisor – Vectorization Optimization

## Faster Vectorization Optimization:

- Vectorize where it will pay off most
- Quickly ID what is blocking vectorization
- Tips for effective vectorization
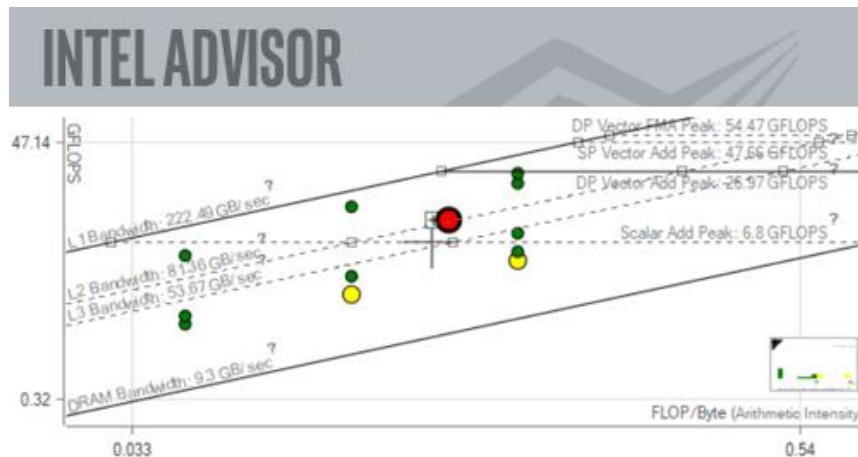- Safely force compiler vectorization
- Optimize memory stride

## Roofline model analysis:

- Automatically generate roofline model
- Evaluate current performance
- Identify boundedness





http://intel.ly/advixe-cl-xe

**Add Parallelism with Less Effort, Less Risk and More Impact**

# Typical Vectorization Optimization Workflow

There is no need to recompile or relink the application, but the use of **-g** is recommended.

Note: if you're using Theta run out of /projects rather than /home

- Collect survey (overhead ~5%) advixe-cl -c survey

  - Basic info (static analysis) - ISA, time spent, etc.

- Collect roofline advixe-cl -c roofline

  - Basically the survey analysis above with roofline analysis (trip counts, flops)

- Collect dependencies (overhead 5-1000x) advixe-cl -c dependencies

  - Differentiate between real and assumed issues blocking vectorization

- Collect Memory Access Patterns advixe-cl -c map

  - Get advice on memory strides

- NB: You can run multiple analyses, and sometimes you have to, to get all of the information you need.

  - For example, in the same batch job you can do a roofline and a dependency analysis and have the output directory be the same so all of that information is shown together in a single context.

  - Make sure you create different output directories for different experiments.

# What is a roofline?

- A roofline is a graphical representation of two factors that affect code performance: flops & memory bandwidth

    - Codes at scale may be limited by file I/O or MPI but there can be different rooflines for those cases.

- Allows you answer questions like:

    - What is/isn't limiting this kernel's performance? Which kernels are more important to overall code performance?

    - What gains might I see from focusing on a particular kernel?

    - Where do I need to focus my software engineering efforts to achieve further gains?

# Roofline cont'd

- Arithmetic intensity (AI)

  - Flops / bytes ratio

    - Bytes can be data moved to/from DRAM, cache, etc.

    - Kernels with a high AI are limited by chip floating point performance (e.g. DGEMM)

    - Kernels with a low AI are limited by memory bandwidth (e.g. STREAM triad, many HPC computational physics kernels)

    - Kernels can be limited by both

# Roofline



Limited by DRAM bandwidth

Limited by floating point ops

$$Gflops = MIN( AI * BW, Gflops_{max})$$

Gflops

Flops/bytes

# Cache-Aware Roofline Optimization
## Next Steps

**If under or near a memory roof…**

- Try a MAP analysis. Make any appropriate **cache optimizations**.
- If cache optimization is impossible, try **reworking the algorithm to have a higher AI.**

**If Under the Vector Add Peak**

Check "Traits" in the Survey to see if FMAs are used. If not, try altering your code or compiler flags to **induce FMA usage.**

**If just above the Scalar Add Peak**

Check **vectorization efficiency** in the Survey. Follow the recommendations to improve it if it's low.

**If under the Scalar Add Peak…**

Check the Survey Report to see if the loop vectorized. If not, try to **get it to vectorize** if possible. This may involve running Dependencies to see if it's safe to force it.



FLOPS

FMA Peak
L1 Bandwidth
Vector Add Peak
L2 Bandwidth
Scalar Add Peak
DRAM Bandwidth

Arithmetic Intensity

# Use --help option!

## advixe-cl --help collect

```
Examples:

Perform a Survey analysis to determine hotspots.

        advisor --collect=survey --project-dir=./advi --search-dir src:r=./src
                -- ./bin/myApplication

Perform a Memory Access Patterns analysis on the specified loops.

        advisor --collect=map --mark-up-list=5,10,12 --project-dir=./advi --search-dir src:r=./src
                -- ./bin/myApplication

Perform a Survey analysis on four nodes of the MPI cluster and store the collected data in the shared ./advi project directory.

        mpirun -n 4 advisor --project-dir=./advi --collect=survey
                -- <PATH>/mpi-sample/1_mpi_sample_serial

Perform a Dependencies analysis on all innermost loops that run above 2% of the total CPU time.

        advisor --collect=dependencies --project-dir=./advi --loops="loop-height=0,total-time>2"
                -- ./bin/myApplication

Perform a Roofline analysis.

        advisor --collect=roofline --project-dir=./advi -- ./bin/myApplication
```

# Generate Advisor Command Lines from the GUI



How accurate you want your reports to be

# Collect survey and tripcounts (roofline)

```
$ Theta: module load PrgEnv-intel amplxe-cl advixe-cl

$ cd project; make

#!/bin/bash

..stuff..

advixe-cl –collect=roofline –trip-counts –project-dir=<project-
dir> -- <executable> <parameters>
```

Where your results go

# View Result on Local Machine

- Make sure your local version of Advisor (or VTune) is at least the same as that of the one used to generate the data otherwise errors might occur.

- X-forwarding is not recommended.

- Tar the result along with sources and binary (if you want to be able to view them, unless you already have them locally)

- Copy to your local machine

- May have to point advixe-cl at your local sources and binary

# Summary Report



Summary provides overall performance characteristics

Top time consuming loops are listed individually

Vectorization efficiency is based on used ISA

# Survey Report (Source Tab)



Notice the following:

- Vector ISA
- Type Conversions
- Memory Access Patterns

All of these elements may affect performance

# Survey Report (Code Analytics Tab)



Analytics tab contains a wealth of information

- Instruction set
- Instruction mix
- Traits (sqrt, type conversions, unpacks)
- Vector efficiency
- Floating point statistics

And explanations on how they are measured or calculated – expand the box or hover over the question marks.

# LIVE DEMO

Roofline

# CARM (Cache-aware roofline model) Analysis



Using single threaded roof

Code vectorized, but performance on par with scalar add peak?

- Irregular memory access patterns force gather operations.

- Overhead of setting up vector operations reduces efficiency.

Next step is clear: perform a Memory Access Pattern analysis

# Memory Access Pattern Analysis (Refinement)

```
advixe-cl –c roofline –r mydat ./nody.x 4000 500

advixe-cl –c map –r mydat ./nbody.x 4000 500
```



Storage of particles is in an Array Of Structures (AOS) style

This leads to regular, but non-unit strides in memory access

- 33% unit

- 0% uniform, non-unit

- 67% non-uniform

Re-structuring the code into a Structure Of Arrays (SOA) may lead to unit stride access and more effective vectorization

# Vectorization: gather/scatter operation

The compiler might generate gather/scatter instructions for loops automatically vectorized where memory locations are not contiguous

```cpp
struct Particle
{
 public:
    ...
    real_type pos[3];
    real_type vel[3];
    real_type acc[3];
    real_type mass;
};
```

```cpp
struct ParticleSoA
{
 public:
    ...
    real_type *pos_x,*pos_y,*pos_z;
    real_type *vel_x,*vel_y,*vel_z;
    real_type *acc_x,*acc_y;*acc_z
    real_type *mass;
};
```

# Memory access pattern analysis

How should I access data ?

## Best: Unit stride access are faster

```
for (i=0; i<N; i++)
    A[i] = B[i]*d
```

## OK: Constant stride are more complex

```
for (i=0; i<N; i+=2)
    A[i] = B[i]*d
```

## Bad: Irregular access

```
for (i=0; i<N; i++)
    A[i] = B[C[i]]*d
```



B

For B, 1 cache line load computes 4 DP

B

For B, 2 cache line loads compute 4 DP with reconstructions

B

For B, 4 cache line loads compute 4 DP with reconstructions, prefetching might not work

# Performance After Data Structure Change

In this new version ( version 3 in GitHub sample ) we introduce the following change:

- Change particle data structures from AOS to SOA

Note changes in report:

- Performance is lower
- Main loop is no longer vectorized
- Assumed vector dependence prevents automatic vectorization



Next step is clear: perform a Dependencies analysis

# Dependencies Analysis (Refinement)

Run "survey" followed by "dependencies"

advixe-cl –c dependencies ./nbody.x 4000 500



Dependencies analysis has high overhead:

- Run on reduced workload

Advisor Findings:

- RAW dependency

# Recommendations

# Performance after resolve dependencies

# Performance After Resolved Dependencies



New memory access pattern plus vectorization produces much improved performance!
What's next? Try suggestions for aligning data.

# Final performance



- Some additional performance eked out.
- Vectorization of loop now 100%
- At this point, you'll likely need to switch to VTune to begin investigating cache misses.

# ADVISOR: GPU OFFLOAD

Which codes to migrate to GPU?

# Offload Advisor

- Another option for accelerating loops is offloading them to an accelerator such as a GPU.

  - As with vectorization, Advisor now has the capability of allowing the user to test if kernels would benefit from offloading

    - Run a number of Advisor collections to generate data

    - Run a projection to a specific architecture (Intel only)

      - Report shows which loops would benefit from offloading, and which would not.

# Nbody test case

- Take final optimized test case (ver7) with the main computational loop parallelized with OpenMP on the host (OpenMP not a requirement)

  - Collect a survey

    - ```
      advixe-cl --collect=survey --project-dir=./advi_proj_v7 --stackwalk-
      mode=online --static-instruction-mix -- ./nbody.x 4000 500
      ```

  - Collect flops and counts and target a particular device

    - ```
      advixe-cl --collect=tripcounts --project-dir=./advi_proj_v7 --flop --
      target-device=gen9_gt2 -- ./nbody.x 4000 500
      ```

  - Do a projection, targeting the same device

    - ```
      advixe-cl --collect=projection  --project-dir=./advi_proj_v7 --
      config=gen9_gt2 --no-assume-dependencies
      ```

# Offload summary

# Offload modeling Accelerated regions tab



- Drill down on offloaded loop
  - Estimated speedup
  - Launch and data transfer latencies
- Offload loop with
  - OpenMP target directives and data mapping clauses
- Profile again with
  - `advixe-cl –collect=roofline –profile-gpu`

# GPU Roofline Summary



- Program time
- GPU time
- Data transfer time
- CPU time
- FPU Utilization
- EU threading Occupancy
- IPC rate
- Thread count
- Roofline
- Hotspots

# GPU Roofline Insights



- Measured GPU roofline
  - L3, SLM, GTI and DRAM bandwidth
- Kernel location
- FPU Utilization
- EU Threading Occupancy
- IPC rate
- Active/stalled/idle %

# Data collection for specific regions: ittnotify

```
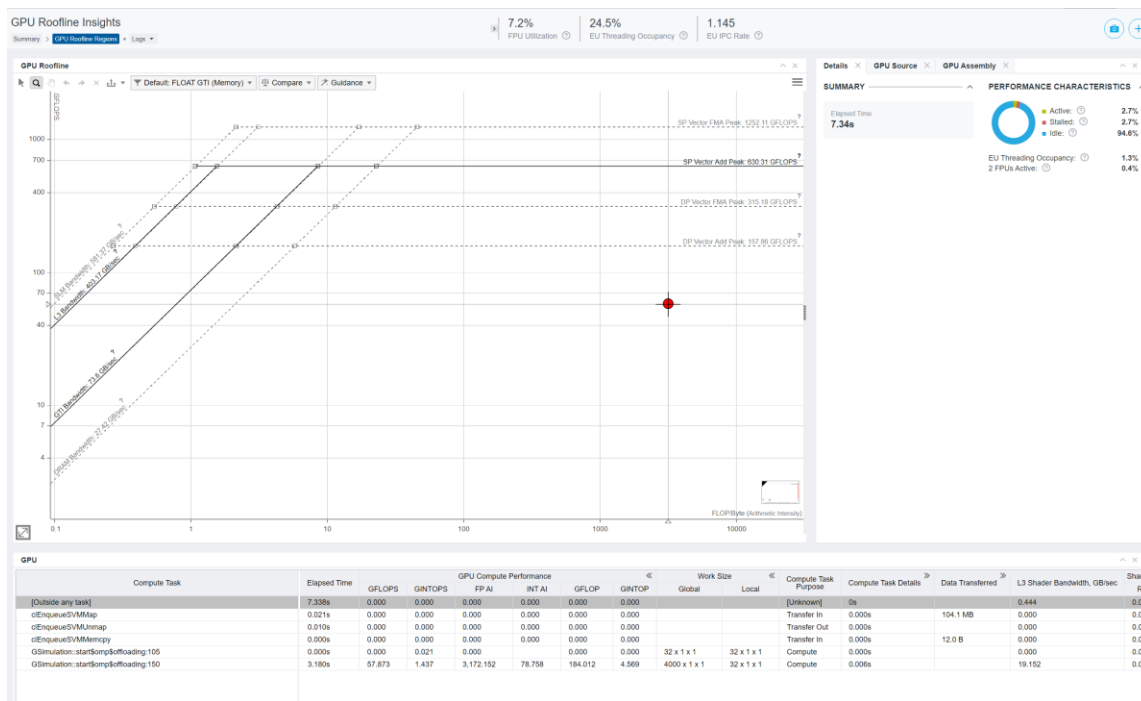#include <ittnotify.h>

Int main(int argc, char* argv[] )

{

// do work here

__itt_pause();

// do more work

__itt_resume()

// Interesting work here

__itt_pause()

// Do more uninteresting work

Return 0;

}
```

- If you just want to collect/examine data from specific regions in your code, you can use ittnotify interface.

  - Link in libittnotify.a

  - Can start program with __itt_pause() or launch with advixe-cl –start-paused.

  - Can also use in VTune

  - Also a Fortran interface

    - e.g. CALL ITT_PAUSE()

# INTEL® VTUNE™

Core-level hardware metrics

https://www.alcf.anl.gov/user-guides/amplxe-cl-xc40

# Intel® VTune™ Amplifier

VTune is a full system (node level) profiler

- Accurate

- Low overhead

- Comprehensive ( CPU, GPU, microarchitecture, memory, IO, threading, ... )

- Configurable interface with easily accessed help

- Direct access to source code and assembly

Analyzing execution behavior with shared resources is critical in achieving good performance on multicore and offload processing systems

# VTune Phases

## Collection

- Record data

- Occurs while target executable is running

## Finalization

- Calculations based on recorded data

- Used in displays / views

- Occurs after collection or in GUI (deferred)

## Reports

- Static report

- Various formats: text, HTML, XML, CSV

## Display (Views)

- VTune GUI

- Interactive, configurable

(intel)

# Predefined Collections

- Many available analysis types (only sme below):

  - hotspots                     Basic hotspots

  - memory-consumption           Use of memory and allocation

  - uarch-exploration            CPU microarchitecture bottlenecks

  - memory-access                Memory access

  - threading                    Threading performance, overhead

  - hpc-performance              OpenMP eff., memory access, vectorization,etc

  - io                           I/O subsystems, CPU, processor buses

  - gpu-offload                  Code execution on cpu and gpu

  - gpu-hotspots                 Hots spots, GPU hw metrics, mem latency, etc

# Copy Command Line to Clipboard

**Command line:**

```
"C:\Program Files (x86)\Intel\oneAPI\vtune\latest\bin64\vtune" -collect hotspots -
app-working-dir "C:\Program Files\DxO\DxO PhotoLab 4" "--app-working-dir=C:\Program
Files\DxO\DxO PhotoLab 4" -- "C:\Program Files\DxO\DxO PhotoLab 4\DxO.PhotoLab.exe"
```

Copy    Close

# HPC-Perf analysis: nbody demo (ver7: threaded)

## Vectorization ⓘ: 100.0% of Packed FP Operations

Instruction Mix:

| | | |
|---|---|---|
| SP FLOPs ⓘ: | 62.3% | of uOps |
|   Packed ⓘ: | 100.0% | from SP FP |
|   Scalar ⓘ: | 0.0% | from SP FP |
| DP FLOPs ⓘ: | 0.0% | of uOps |
|   Packed ⓘ: | 0.0% | from DP FP |
|   Scalar ⓘ: | 0.0% | from DP FP |
| x87 FLOPs ⓘ: | 0.0% | of uOps |
| Non-FP ⓘ: | 37.7% | of uOps |
| FP Arith/Mem Rd Instr. Ratio ⓘ: 6.662 | | |
| FP Arith/Mem Wr Instr. Ratio ⓘ: 495.224 | | |

*N/A is applied to metrics with undefined value. There is no data to calculate the metric.

### Top Loops/Functions with FPU Usage by CPU Time
This section provides information for the most time consuming loops/functions with floating point operations.

| Function | CPU Time ⓘ | % of FP Ops ⓘ | FP Ops: Packed ⓘ | FP Ops: Scalar ⓘ | Vector Instruction Set ⓘ | Loop Type ⓘ |
|---|---|---|---|---|---|---|
| [Loop at line 193 in GSimulation::start$omp$parallel@163] | 8.234s | 62.9% | 100.0% | 0.0% | AVX(256); FMA(256) | Body |
| [Loop at line 165 in GSimulation::start$omp$parallel@163] | 0.023s | 25.0% | 100.0% | 0.0% | AVX(128); AVX(256); AVX2(256); FMA(256) | Body |

*N/A is applied to non-summable metrics.

# HPC-Perf: Bottom-up Hotspots view



**Use drop down menu to access 'Hotspots by CPU Utilization'**

Double click on line to access source and assembly.

Notice the filtering options at the bottom, which allow customization of this view.

Can also do this under "HPC Performance Characterization" and see loop/function data for spin time, serialization, FP Ops, CPI, etc.

Next steps would include additional analysis to continue the optimization process.

# HPC Perf: Bottoms Up – Source View



Click through bottom's up view to see source and metrics.

# HPC Perf: Memory Usage



## Get overview of

- total loads/stores

- bandwidth usage

- L3 bandwidth

- GPU bandwidths

- top functions with high bandwidth utilization.

# HPC Perf: Memory Usage, bottom's up view



Can see loads and stores by loop/function

- Sort by loads/stores/llc miss counts

Can also click through as before to see source level view of

- Cpu time

- Loads/stores

- LLC miss counts

# HPC-Perf: Hardware Events



**Drop down for Hardware Events**.

Summary of all measured performance counters

Very similar results to what you'd get from 'uarch-collection'

# HPC-Perf: Hardware Performance Counters



- **Bottom up view of counters**

- Scroll window to see all counters.

- Timeline of counter activity

- Click through to see source Level view of counter data

- CPU thread and GPU counters

# uarch-exploration: summary



- Running 4 threads

- Want 'Retiring to be 100%' – high instruction throughput

- Reporting core bound – implying not enough resources available, in this case likely to be FP units

- Also can look at raw performance counters (including timeline).

# uarch-exploration: bottom up



Like other collections, can click through the top hotspots to see source code and where limiter is seen to be sequence of operations with high flop counts.

# uarch-exploration: 8 thread summary



- Running with 8 threads improves performance very slightly but shows code is now frontend bound

  - Likely due to pipeline slots being stalled due to too many memory references per cycle.

# Memory-Consumption collection:



Shows  top memory consumers

Bottom's up show's by loop/function/timeline of consumption

Can click through function/loop to see allocation/deallocation sizes at source level.

# VTUNE: ITT

# Itt pause & resume

Use `__itt_pause()` &
`__itt_resume()` to target
data collection only in
specific regions.

```
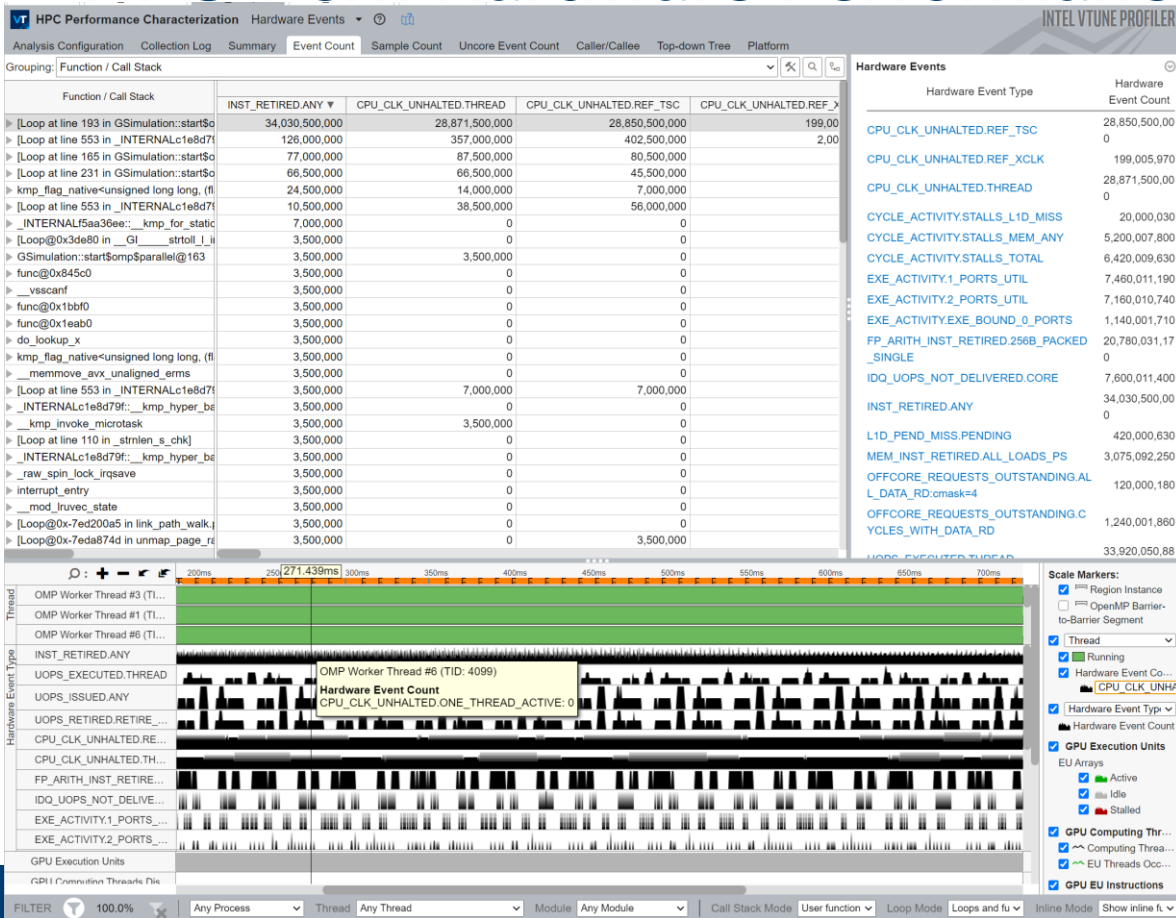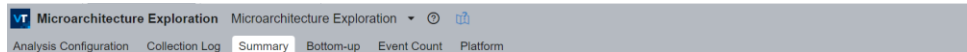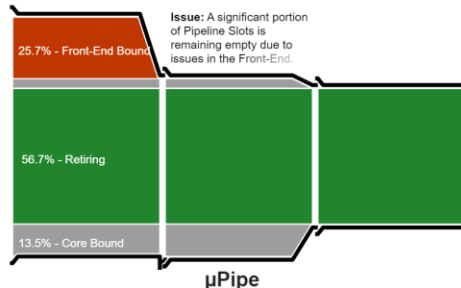#include <ittnotify.h>
    ...uninteresting work....
    __itt_resume();
    ...interesting work....
    __itt_pause();
    ...more uninteresting work...
```

- Launch with `amplxe-cl –start-paused`
  ......

# Using itt to create custom counters

One can create custom counters that show up on VTune timelines by using the itt interface.

In the example at the right, the counter "myFlops" will show up in the performance metrics timelines.

Tested with hotspot and uarch-exploration and it works. Some issue with hpc-performance that is being looked at.

```
#include <ittnotify.h>

Main()
{

    __itt_counter  myCounter;
    __itt_counter_create("myFlops", "Domain");
    …do some stuff….
    __itt_counter__set_value( myCounter, &val );
    …do some stuff…
    __itt_counter_set_value( myCounter, &val);
    …do more stuff…
    __itt_counter_inc_delta( myCounter, &val );
    …do more stuff….
    __itt_counter_dec_delta( myCounter, &val);

    __itt_counter_destroy(myCounter);

}
```

# Counter creation example

# Collected using 'hotspots'

# VTUNE: TARGETING MPI RANKS

# Collecting on Single MPI Ranks

- Might want to use VTune on an MPI application but not, by default, collect data on all MPI ranks as VTune is not designed for that.

  - Still possible to gather some useful data.

  - Using ittnotify is not the route as it still collects data on all ranks even if you pause collection before MPI_Init()

  - Use env vars and MPMD mode:

```
mpirun –genv I_MPI_PIN_PROCESSOR_LIST=0-2,4-7 –n 7 ./app :
        -genv I_MPI_PIN_PROCESSOR_LIST=3 –n 1 amplxe-cl –c
hotspots
        –r appdat -- ./app
```

# VTUNE: GPU OFFLOADING

# VTune gpu-offloading

- Graphics Information about speeds and feeds (no context)

# PROFILING PYTHON & ML APPLICATIONS

# Python

Profiling Python is straightforward in VTune™ Amplifier, as long as one does the following:

- The "application" should be the full path to the python interpreter used

- The python code should be passed as "arguments" to the "application"

In Theta this would look like this:

```
mpirun -n 1 -N 1 amplxe-cl -c hotspots -r res_dir \
                -- /usr/bin/python3 mycode.py myarguments
```

# Simple Python Example on Theta

```
mpirun -n 1 -N 1 amplxe-cl -c hotspots -r vt_pytest \
                    -- /usr/bin/python ./cov.py naive 100 1000
```



Naïve implementation of the calculation of a covariance matrix

Summary shows:

- Single thread execution

- Top function is "naive"

Click on top function to go to Bottom-up view

# Bottom-up View and Source Code



Inefficient array multiplication found quickly
We could use numpy to improve on this

Note that for mixed Python/C code a Top-Down view can often be helpful to drill down into the C kernels

# COMMON ISSUES

# Fixes

No call stack information/unknown stack frame

- Check finalization log
    - Make sure VTune finds your binary along with libraries that you call

Incompatible database scheme when trying to open result in GUI

- Make sure your local VTune is the same version or newer

VTune sampling driver.. using perf or errors mentioning PMU Resources

- Notify support@alcf.anl.gov  or your nearest Intel COE person

# TIPS AND TRICKS

# Speeding up finalization

**Advisor**

add `--no-auto-finalize` to the aprun

followed by `advixe-cl R survey ...` _without aprun_ will cause to finalize on the momnode rather than KNL.

You can also finalize on thetalogin:

cd your_src_dir;

export SRCDIR=`pwd | xargs realpath`

advixe-cl -R survey --search-dir src:=${SRCDIR} ..

**VTune**

add `--finalization-mode=none` to aprun

followed by `amplxe-cl -R hotspots ...` _without aprun_ will cause to finalize on momnode rather than KNL

You can also finalize on thetalogin:

cd your_src_dir;

export SRCDIR=`pwd | xargs realpath`

amplxe-cl -R hotspots --search-dir src:=${SRCDIR} ..

# Managing overheads

Advisor Dependencies and MAP analyses can have huge overheads

If able, run on reduced problem size. Advisor just needs to figure out the execution flow.

Only analyze loops/functions of interest:

> https://software.intel.com/en-us/advixe-cl-user-guide-mark-up-loops

# When do I use VTune vs Advisor?

## VTune

- What's my cache hit ratio?

- Which loop/function is consuming most time overall? (bottom-up)

- Am I stalling often? IPC?

- Am I keeping all the threads busy?

- Am I hitting remote NUMA?

- When do I maximize my BW?

## Advisor

- Which vector ISA am I using?

- Flow of execution (callstacks)

- What is my vectorization efficiency?

- Can I safely force vectorization?

- Inlining? Data type conversions?

- Roofline

# BACKUP

# VTune Cheat Sheet

Compile with –g -dynamic

amplxe-cl –c hpc-performance –flags -- ./executable

- --result-dir=./amplxe-cl_output_dir

- --search-dir src:=../src --search-dir bin:=./

- -knob enable-stack-collection=true –knob collect-memory-bandwidth=false

- -knob analyze-openmp=true

- -finalization-mode=deferred if finalization is taking too long on KNL

- -data-limit=125 ← in mb

- -trace-mpi for MPI metrics on Theta

- amplxe-cl –help collect survey

**https://software.intel.com/en-us/amplxe-cl-amplifier-help-amplxe-cl-command-syntax**

(intel)

# Advisor Cheat Sheet

Compile with –g -dynamic

advixe-cl –c roofline/depencies/map –flags -- ./executable

- --project-dir=./advixe_output_dir

- --search-dir src:=../src --search-dir bin:=./

- -no-auto-finalize if finalization is taking too long on KNL

- --interval 1 (sample at 1ms interval, helps for profiling short runs)

- -data-limit=125 ← in mb

- advixe-cl -help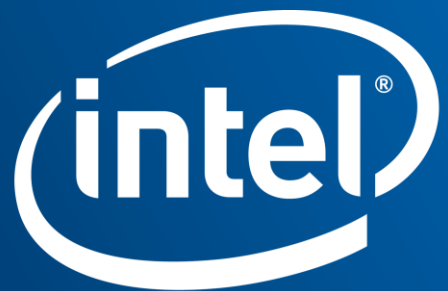